# 9th USENIX Symposium on Operating Systems Design and Implementation (OSDI'10)

Vancouver, Canada
4 - 6 October 2010

# 9th USENIX Symposium on Operating Systems Design and Implementation
## October 4–6, 2010
## Vancouver, BC, Canada

## Monday, October 4

### Kernels: Past, Present, and Future

### Inside the Data Center, 1

### Security Technologies

### Concurrency Bugs

## Tuesday, October 5

## Wednesday, October 6